
Conception d'un banc d'essais décisionnel

Jérôme Darmont — Fadila Bentayeb — Omar Boussaïd

ERIC – Université Lumière Lyon 2
5 avenue Pierre Mendès-France
69676 Bron Cedex
Contact : jerome.darmont@univ-lyon2.fr

RÉSUMÉ. Nous présentons dans cet article un nouveau banc d'essais pour l'évaluation des performances des entrepôts de données. L'emploi de bancs d'essais est profitable pour les utilisateurs, qui peuvent comparer les performances de plusieurs systèmes avant d'investir, mais aussi pour les concepteurs d'entrepôts de données, afin d'évaluer l'impact de différents choix techniques (indexation, matérialisation de vues...). Les bancs d'essais standards édités par le TPC (Transaction Processing Performance Council) répondent au premier de ces besoins, mais ne sont pas suffisamment adaptables pour satisfaire le second. C'est pourquoi nous proposons le banc d'essais DWEB (Data Warehouse Engineering Benchmark), qui permet de générer à la demande divers entrepôts de données synthétiques, ainsi que les charges (ensembles de requêtes décisionnelles) associées. DWEB est totalement adaptable, mais il demeure facile à mettre en œuvre grâce à deux niveaux de paramétrage. Par ailleurs, comme DWEB répond principalement à des besoins d'évaluation de performance pour l'ingénierie, il est complémentaire plutôt que concurrent aux bancs d'essais standards du TPC. Finalement, DWEB est implémenté sous la forme d'un logiciel libre écrit en Java qui peut s'interfacer avec la plupart des systèmes de gestion de bases de données relationnels existants.

ABSTRACT. We present in this paper a new benchmark for evaluating the performances of data warehouses. Benchmarking is useful either to system users for comparing the performances of different systems, or to system engineers for testing the effect of various design choices. While the TPC (Transaction Processing Performance Council) standard benchmarks address the first point, they are not tuneable enough to address the second one. Our Data Warehouse Engineering Benchmark (DWEB) allows to generate various ad-hoc synthetic data warehouses and workloads. DWEB is fully parameterized. However, two levels of parameterization keep it easy to tune. Since DWEB mainly meets engineering benchmarking needs, it is complimentary to the TPC standard benchmarks, and not a competitor. Finally, DWEB is implemented as a Java free software that can be interfaced with most existing relational database management systems.

MOTS-CLÉS : Entrepôts de données, requêtes décisionnelles, OLAP, bancs d'essais, évaluation de performance, conception d'entrepôts de données.

KEYWORDS : Data warehouses, decision support queries, OLAP, benchmarking, performance evaluation, data warehouse design.

1. Introduction

Evaluer des technologies centrées sur la décision telles que les entrepôts de données n'est pas une tâche facile. Bien que des conseils pertinents d'ordre général soient disponibles en ligne [PEN 03, GRE 04a], les éléments plus quantitatifs au regard des performances de ces systèmes sont rares.

Dans le contexte des bases de données transactionnelles, des bancs d'essais sont traditionnellement utilisés pour évaluer la performance. Ce sont des modèles synthétiques de bases de données et de charges (opérations à effectuer sur la base), ainsi que des ensembles de mesures de performance. Dans le contexte de l'aide à la décision et plus précisément lors de la conception et de l'exploitation d'un entrepôt de données, de bonnes performances sont encore plus critiques en raison de la nature du modèle de données spécifique et de la volumétrie de ces données. L'objectif de cet article est de proposer un nouveau banc d'essais pour entrepôts de données.

Plusieurs objectifs peuvent être visés lors de l'utilisation d'un banc d'essais :

- 1) comparer les performances de plusieurs systèmes dans des conditions expérimentales données (utilisateurs de ces systèmes) ;
- 2) évaluer l'impact de différents choix architecturaux ou de techniques d'optimisation sur les performances d'un système donné (concepteurs d'entrepôts de données).

Les bancs d'essais proposés par le TPC (*Transaction Processing Performance Council*), un organisme à but non lucratif qui définit des bancs d'essais standards et publie les résultats d'évaluations de performance de façon indépendante, répondent parfaitement au premier de ces objectifs. Cependant, ils ne sont pas très adaptables : leur seul paramètre est un facteur qui définit la taille globale de leur base de données. Néanmoins, dans un contexte de développement, il est intéressant de pouvoir tester une solution (une stratégie d'indexation automatique, par exemple) sur différentes configurations de base de données. Notre objectif est donc de proposer un banc d'essais permettant de générer des entrepôts de données synthétiques, ainsi que les charges (ensembles de requêtes décisionnelles) associées, pour satisfaire des besoins d'ingénierie. Nous l'avons baptisé DWEB (*Data Warehouse Engineering Benchmark*). Il faut souligner que DWEB n'est pas concurrent des bancs d'essais décisionnels standards édités par le TPC. En effet, nous le considérons plutôt comme complémentaire, puisqu'il cible principalement le second objectif mentionné plus haut.

Cet article est organisé comme suit. Nous étudions tout d'abord l'état de l'art concernant les bancs d'essais décisionnels dans la section 2. Nous détaillons ensuite la base de données et la charge de DWEB dans les sections 3 et 4, respectivement. Nous présentons également brièvement notre implémentation de DWEB dans la section 5. Nous concluons finalement et présentons nos perspectives de recherche dans la section 6.

2. Bancs d'essais décisionnels existants

À notre connaissance, il existe très peu de bancs d'essais décisionnels en-dehors de ceux du TPC. De plus, les spécifications de ceux que nous avons recensés sont rarement publiées dans leur intégralité [DEM 95]. C'est pourquoi nous nous concentrons sur les bancs d'essais du TPC dans cette section.

TPC-D [BAL 93, BHA 96, Tra98] a fait son apparition dans le milieu des années 90 et forme la base de TPC-H et TPC-R [POE 00, Tra03a, Tra03b], qui l'ont désormais remplacé. TPC-H et TPC-R sont en fait identiques, seul leur mode d'utilisation les différencie. TPC-H est conçu pour le requêtage ad-hoc (les requêtes ne sont pas connues à l'avance et toute optimisation préalable est interdite), tandis que TPC-R a une vocation de *reporting* (les requêtes sont connues à l'avance et des optimisations adéquates sont possibles). TPC-H et TPC-R exploitent le même schéma de base de données que TPC-D : un modèle *produits-commandes-fournisseurs* classique (représenté par un diagramme de classes UML dans la figure 1); ainsi que la charge de TPC-D enrichie de cinq nouvelles requêtes.

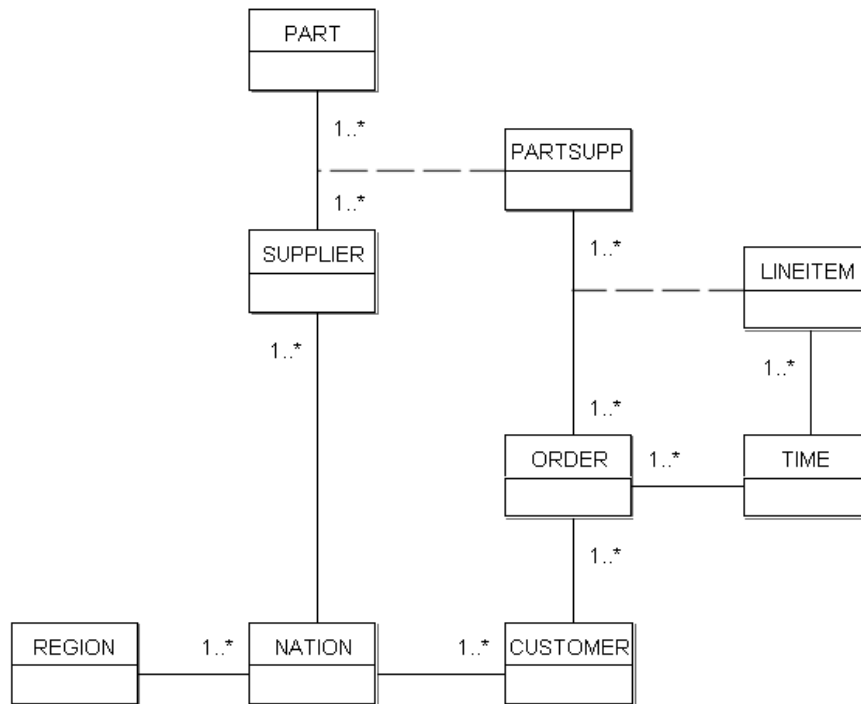


Figure 1 – Schéma de la base de données de TPC-D, TPC-H et TPC-R

Plus précisément, cette charge est constituée de vingt-deux requêtes décisionnelles

paramétrées écrites en SQL-92 et numérotées Q1 à Q22 et de deux fonctions de rafraîchissement RF1 et RF2 qui ajoutent et suppriment des n-uplets dans les tables ORDER et LINEITEM. Les paramètres des requêtes sont instanciés aléatoirement en suivant une loi uniforme. Finalement, le protocole d'exécution de TPC-H ou TPC-R est le suivant :

- 1) un test de chargement ;
- 2) un test de performance (exécuté deux fois), lui même subdivisé en un test de puissance et un test de débit.

Trois mesures principales permettent de décrire les résultats obtenus en termes de puissance, de débit et d'une composition de ces deux critères.

TPC-DS, qui est actuellement en cours de développement [POE 02], modélise plus clairement un entrepôt de données. Il est le successeur annoncé de TPC-H et TPC-R. Le schéma de la base de données de TPC-DS, dont les tables de faits sont représentées dans la figure 2, représente les fonctions décisionnelles d'un détaillant sous la forme de plusieurs schémas en flocon de neige. Ce modèle comprend également quinze dimensions partagées par les tables de faits. Il s'agit donc d'un schéma en constellation.

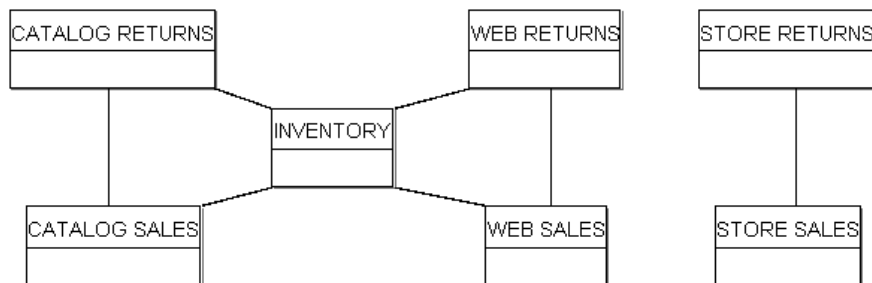


Figure 2 – Schéma de l'entrepôt de données de TPC-DS (tables de faits)

La charge de TPC-DS est constituée de quatre classes de requêtes : requêtes de *reporting*, requêtes décisionnelles ad-hoc, requêtes interactives d'analyse en ligne (OLAP – *On-Line Analytical Processing*), requêtes d'extraction. Des modèles de requêtes écrits en SQL-99 (et comprenant donc des extensions OLAP) permettent de générer un ensemble d'environ cinq cents requêtes. Ces modèles sont instanciés aléatoirement selon des distributions non-uniformes. Le processus de maintenance de l'entrepôt de données comprend une phase d'ETL (*Extract, Transform, Load*) complète et un traitement spécifique des dimensions. Par exemple, les dimensions historisées conservent les anciens n-uplets quand de nouveaux sont ajoutés, tandis que les dimensions non-historisées ne conservent pas les anciennes données. Finalement, le modèle d'exécution de TPC-DS est divisé en quatre étapes :

- 1) un test de chargement,
- 2) une exécution des requêtes,

- 3) une phase de maintenance des données,
- 4) une seconde exécution des requêtes.

Une seule mesure (de débit) est proposée. Elle prend en compte l'exécution des requêtes et la phase de maintenance.

Bien que les bancs d'essais décisionnels du TPC soient adaptables selon la définition de Gray [GRA 93], leur schéma est fixe et ils ne sont pas idéaux pour évaluer l'impact de choix architecturaux ou de techniques d'optimisation sur les performances globales. En effet, un seul paramètre permet de définir leur base de données (*Scale Factor – SF*) en terme de taille (de 1 à 100 000 Go). De plus, leur charge n'est pas du tout paramétrable : le nombre de requêtes générées dépend directement de *SF* dans TPC-DS, par exemple. Il s'avère donc intéressant de proposer un banc d'essais complémentaire capable de modéliser diverses configurations d'entrepôts de données.

3. Base de données de DWEB

3.1. Schéma

Notre objectif avec DWEB est de pouvoir modéliser différents types d'architectures d'entrepôts de données [INM 02, KIM 02] au sein d'un environnement ROLAP (*Relational OLAP*) :

- des schémas en étoile classiques,
- des schémas en flocon de neige (avec des dimensions hiérarchisées),
- des schémas en constellation (avec plusieurs tables de faits et des dimensions partagées).

Afin d'atteindre ce but, nous proposons un métamodèle d'entrepôt de données (représenté par un diagramme de classes UML dans la figure 3) susceptible d'être instancié en ces différents schémas. Nous considérons ce métamodèle comme un intermédiaire entre le métamodèle multidimensionnel du standard CWM (*Common Warehouse Metamodel*) [Obj03, POO 03] et le modèle final du banc d'essais. Notre métamodèle est en fait une instance de celui de CWM, qui pourrait en fait être qualifié de méta-métamodèle dans notre contexte.

Notre métamodèle est relativement simple. La partie supérieure de la figure 3 décrit un entrepôt de données (ou un magasin de données si ce dernier est perçu comme un petit entrepôt dédié) constitué d'une ou plusieurs tables de faits qui sont chacune décrites par plusieurs dimensions. Chaque dimension peut également décrire plusieurs tables de faits (dimensions partagées) et peut contenir une ou plusieurs hiérarchies composées de différents niveaux. Il est possible de n'avoir qu'un seul niveau, auquel cas la dimension n'est pas hiérarchisée.

Les tables de faits et les niveaux hiérarchiques des dimensions sont tous des tables relationnelles, qui sont modélisées dans la partie inférieure de la figure 3. Classique-

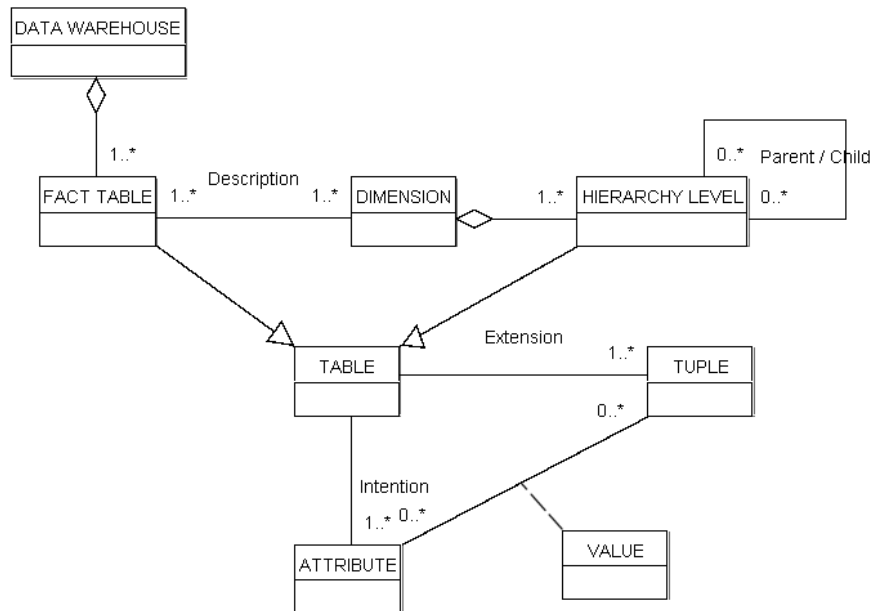


Figure 3 – Métaschéma de l'entrepôt de données de DWEB

ment, une table ou relation est définie en intention par ses attributs et en extension par ses n-uplets. À l'intersection d'un attribut et d'un n-uplet donnés se trouve la valeur de cet attribut dans ce n-uplet.

3.2. Paramétrage

La difficulté principale dans le processus de création d'un schéma d'entrepôt de données est le paramétrage de l'instanciation du métaschéma. Nous avons en effet pour objectif de satisfaire les quatre critères qui définissent, d'après Gray, un « bon » banc d'essais [GRA 93] :

- *pertinence* : le banc d'essais doit répondre aux besoins d'ingénierie exprimés dans la section 1 ;
- *portabilité* : le banc d'essais doit être facile à implémenter sur différents systèmes ;
- *adaptabilité* : il doit être possible d'étudier des bases de données de tailles diverses et d'augmenter l'échelle du banc d'essais ;
- *simplicité* : le banc d'essais doit être compréhensible sous peine de ne pas être crédible et de demeurer inutilisé.

La pertinence et la simplicité forment clairement des objectifs orthogonaux. Introduire des paramètres en trop petit nombre réduit l'expressivité du modèle, tandis qu'en prévoir trop rend le banc d'essais difficile à appréhender par des utilisateurs potentiels. En parallèle, la complexité de génération du schéma instancié doit être maîtrisée. Pour résoudre ce dilemme, nous tirons partie de l'expérience de la conception du banc d'essais orienté objets OCB [DAR 00]. OCB est générique et capable de simuler tous les autres bancs d'essais orientés objets, mais sa mise en œuvre est contrôlée par des paramètres trop nombreux, si bien que seule une minorité d'entre eux est utilisée en pratique. Nous proposons donc de subdiviser les paramètres en deux sous-ensembles :

- un sous-ensemble détaillé de paramètres de bas niveau qui permet à un utilisateur avancé de contrôler totalement la génération de l'entrepôt de données (tableau 1) ;
- une couche supérieure contenant beaucoup moins de paramètres facilement compréhensibles et instanciables (tableau 2). Ces paramètres de haut niveau sont plus précisément les valeurs moyennes des paramètres de bas niveau. Lors de la génération de la base de données, ils sont exploités par des fonctions aléatoires (qui suivent une distribution gaussienne) pour fixer la valeur des paramètres de bas niveau. Finalement, bien que le nombre de paramètres de bas niveau puisse augmenter de façon significative si le schéma grossit, le nombre de paramètres de haut niveau, lui, demeure constant et raisonnable (moins de dix paramètres).

Les utilisateurs peuvent alors choisir d'initialiser l'ensemble complet de paramètres de bas niveau ou bien seulement les paramètres de haut niveau, pour lesquels nous proposons des valeurs par défaut correspondant à un schéma en flocon de neige.

Nom du paramètre	Signification
NB_FT	Nombre de tables de faits
$NB_DIM(f)$	Nombre de dimensions décrivant la table de faits f
TOT_NB_DIM	Nombre total de dimensions
$NB_MEAS(f)$	Nombre de mesures dans la table de faits f
$DENSITY(f)$	Densité de la table de faits f
$NB_LEVELS(d)$	Nombre de niveaux hiérarchiques dans la dimension d
$NB_ATT(d, h)$	Nombre d'attributs dans le niveau hiérarchique h de la dimension d
$HHLEVEL_SIZE(d)$	Nombre de n-uplets dans le plus haut niveau hiérarchique de la dimension d
$DIM_SFACTOR(d)$	Facteur d'échelle pour la taille des niveaux hiérarchiques de la dimension d

Tableau 1 – Paramètres de bas niveau de l'entrepôt de données de DWEB

Remarques :

- Puisque des dimensions partagées sont possibles,

$$TOT_NB_DIM \leq \sum_{i=1}^{NB_FT} NB_DIM(i).$$

Nom du paramètre	Signification	Valeur déf.
<i>AVG_NB_FT</i>	Nombre moyen de tables de faits	1
<i>AVG_NB_DIM</i>	Nombre moyen de dimensions par table de fait	5
<i>AVG_TOT_NB_DIM</i>	Nombre moyen total de dimensions	5
<i>AVG_NB_MEAS</i>	Nombre moyen de mesures dans les tables de faits	5
<i>AVG_DENSITY</i>	Densité moyenne des tables de faits	0,6
<i>AVG_NB_LEVELS</i>	Nombre moyen de niveaux hiérarchiques dans les dimensions	3
<i>AVG_NB_ATT</i>	Nombre moyen d'attributs dans les niveaux hiérarchiques	5
<i>AVG_HHLEVEL_SIZE</i>	Nombre moyen de n-uplets dans les plus hauts niveaux hiérarchiques	10
<i>DIM_SFCTOR</i>	Facteur d'échelle de taille moyen au sein des niveaux hiérarchiques	10

Tableau 2 – Paramètres de haut niveau de l'entrepôt de données de DWEB

– La cardinalité d'une table de faits est habituellement inférieure ou égale au produit des cardinalités de ses dimensions. C'est pourquoi nous introduisons la notion de densité. Une densité de 1 indique que toutes les combinaisons possibles des clés primaires des dimensions sont présentes dans la table de faits. Quand la densité diminue, nous éliminons progressivement certaines de ces combinaisons (voir section 3.3).

– Au sein d'une dimension, un niveau hiérarchique donné a normalement une cardinalité plus grande que le niveau suivant. Par exemple, dans une hiérarchie de type *villes-régions-pays*, le nombre de villes doit être plus grand que le nombre de régions, qui doit à son tour être supérieur au nombre de pays. De plus, il existe souvent un facteur d'échelle significatif entre ces cardinalités (par exemple, mille villes, cent régions, dix pays). C'est pourquoi nous définissons la cardinalité des niveaux hiérarchiques en assignant une cardinalité « de départ » au plus haut niveau de la hiérarchie (*HHLEVEL_SIZE*). Nous la multiplions ensuite par un facteur d'échelle prédéfini (*DIM_SFCTOR*) pour chaque niveau hiérarchique inférieur.

3.3. Algorithme de génération

L'instanciation du métaschéma de DWEB en un schéma d'entrepôt de données réel s'effectue en deux étapes :

- 1) construction des dimensions et de leurs hiérarchies ;
- 2) construction des tables de faits.

L'algorithme correspondant à ces deux étapes est présenté dans les figures 4 et 5, respectivement. Chacune d'entre elles est constituée, pour chaque dimension ou table

de fait, de la génération de son intention, puis de son extension. De plus, les hiérarchies des dimensions doivent être gérées. Il faut noter qu'elles sont générées en démarrant du plus haut niveau hiérarchique. Par exemple, pour notre hiérarchie *villes-régions-pays*, nous construisons le niveau *pays* en premier, puis le niveau *région* et enfin le niveau *ville*. Ainsi, les n-uplets d'un niveau hiérarchique donné peuvent référencer ceux du niveau supérieur (qui sont déjà créés) à l'aide d'une clé étrangère.

```

For i = 1 to TOT_NB_DIM do
  previous_ptr = NIL
  size = HHLEVEL_SIZE(i)
  For j = 1 to NB_LEVELS(i) do
    // Intention
    hl = New(Hierarchy_level)
    hl.intention = Primary_key()
    For k = 1 to NB_ATT(i,j) do
      hl.intention = hl.intention ∪ String_descriptor()
    End for
    // Gestion des hiérarchies
    hl.child = previous_ptr
    hl.parent = NIL
    If previous_ptr ≠ NIL then
      previous_ptr.parent = hl
      hl.intention = hl.intention
      ∪ previous_ptr.intention.primary_key // Clé étrangère
    End if
    // Extension
    hl.extension = ∅
    For k = 1 to size do
      new_tuple = Integer_primary_key()
      For l = 1 to NB_ATT(i,j) do
        new_tuple = new_tuple ∪ Random_string()
      End for
      If previous_ptr ≠ NIL then
        new_tuple = new_tuple ∪ Random_key(previous_ptr)
      End if
      hl.extension = hl.extension ∪ new_tuple
    End for
    previous_ptr = hl
    size = size * DIM_SFACTOR(i)
  End for
  dim(i) = hl // Premier (plus bas) niveau de la hiérarchie
End for

```

Figure 4 – Algorithme de génération des dimensions

Nous utilisons trois classes principales de fonctions et une procédure dans ces algorithmes.

```

For i = 1 to NB_FT do
  // Intention
  ft(i).intention = ∅
  For j = 1 to NB_DIM(i) do
    j = Random_dimension(ft(i))
    ft(i).intention = ft(i).intention ∪ ft(i).dim(j).primary_key
  End for
  For j = 1 to NB_MEAS(i) do
    ft(i).intention = ft(i).intention ∪ Float_measure()
  End for
  // Extension
  ft(i).extension = ∅
  For j = 1 to NB_DIM(i) do // Produit cartésien
    ft(i).extension = ft(i).extension × ft(i).dim(j).primary_key
  End for
  to_delete = DENSITY(i) * |ft(i).extension|
  For j = 1 to to_delete do
    Random_delete(ft(i).extension)
  End for
  For j = 1 to |ft(i).extension| do
    // |ft(i).extension| est supposée à jour
    For k = 1 to NB_MEAS(i) do
      ft(i).extension.tuple(j).measure(k) = Random_float()
    End for
  End for
End for

```

Figure 5 – Algorithme de génération des tables de faits

1) `Primary_key()`, `String_descriptor()` et `Float_measure()` renvoient des noms pour les clés primaires, les descripteurs des niveaux hiérarchiques et les mesures des tables de faits, respectivement. Ces noms sont étiquetés séquentiellement et préfixés par le nom de la table (par exemple, `DIM1_1_DESCR1`, `DIM1_1_DESCR2...`).

2) `Integer_primary_key()`, `Random_key()`, `Random_string()` et `Random_float()` renvoient des entiers séquentiels par rapport à une table donnée (aucun doublon n'est permis), des instances aléatoires de la clé primaire de la table spécifiée (valeurs aléatoires pour une clé étrangère), des chaînes de caractères aléatoires de taille fixe (20 caractères) sélectionnées dans un référentiel préalablement construit et préfixées par le nom de l'attribut correspondant, ainsi que des nombres aléatoires réels simple précision, respectivement.

3) `Random_dimension()` renvoie une dimension sélectionnée parmi les dimensions existantes qui ne décrivent pas déjà la table de faits passée en argument.

4) `Random_delete()` efface aléatoirement un n-uplet dans l'extension d'une table.

A l'exception de la procédure `Random_delete()`, dans laquelle nous utilisons une distribution aléatoire uniforme, nous employons des distributions aléatoires gaussiennes.

Remarque : La manière dont la densité est gérée dans la figure 5 est clairement non-optimale. Nous choisissons de présenter l'algorithme de cette manière afin de le rendre plus clair, mais notre implémentation ne crée pas tous les n-uplets du résultat du produit cartésien avant d'en effacer certains. Nous générons directement le bon nombre de n-uplets en utilisant la densité comme une probabilité de création de chaque n-uplet.

4. Charge de DWEB

Dans un banc d'essais pour entrepôts de données, la charge peut être subdivisée en :

- une charge de requêtes décisionnelles (principalement des requêtes OLAP) ;
- le processus d'ETL (chargement et rafraîchissement des données).

TPC-DS gère ces deux aspects. Puisque la base de données de DWEB compose la partie la plus originale de notre travail, nous nous inspirons donc de la définition de charge de TPC-DS, tout en exploitant d'autres sources d'informations concernant la performance des entrepôts de données [BMC 00, GRE 04b]. Néanmoins, ils nous est nécessaire d'adapter les propositions de TPC-DS à la nature variable des entrepôts de données de DWEB. De plus, nous voulons satisfaire le critère de simplicité de Gray, aussi proposons-nous au final une charge plus simple que celle de TPC-DS.

Par ailleurs, nous nous concentrons dans un premier temps sur la définition d'un modèle de requêtes. La modélisation du processus d'ETL complet est une tâche complexe sur laquelle nous comptons revenir plus tard. Nous considérons pour l'instant que les spécifications actuelles de DWEB permettent une évaluation grossière d'un chargement d'entrepôt. En effet, la base de données de DWEB peut être générée sous forme de fichiers plats, puis ensuite chargée dans un entrepôt en utilisant les outils fournis par le système.

4.1. Modèle de requêtes

La charge de DWEB modélise deux types de requêtes :

- des requêtes purement décisionnelles contenant les opérations OLAP usuelles telles que cube, agrégation (*roll-up*), forage (*drill down*) et projection/sélection (*slice and dice*) ;
- des requêtes d'extraction.

Nous définissons notre modèle de requêtes générique (figure 6) à l'aide d'une grammaire qui est un sous-ensemble du standard SQL-99, qui introduit des outils ana-

lytiques (indispensables dans un contexte décisionnel) dans les requêtes relationnelles. Cela nous permet de générer des requêtes SQL analytiques dynamiquement.

Query :-	
Select	! [<Attribute Clause> <Aggregate Clause> [<Attribute Clause>, <Aggregate Clause>]]
From	! <Table Clause> [<Where Clause> [<Group by Clause> * <Having Clause>]]
Attribute Clause :-	Attribute Name [[, <Attribute Clause>] \perp]
Aggregate Clause :-	! [Aggregate Function Name (Attribute Name)] [As Alias [[, <Aggregate Clause>] \perp]
Table Clause :-	Table Name [[, <Table Clause>] \perp]
Where Clause :-	Where ! [<Condition Clause> <Join Clause> [<Condition Clause> And <Join Clause>]]
Condition Clause :-	! [Attribute Name <Comparison Operator> <Operand Clause>] [[<Logical Operator> <Condition Clause>] \perp]
Operand Clause :-	[Attribute Name Attribute Value Attribute Value List]
Join Clause :-	! [Attribute Name i = Attribute Name j] [[And <Join Clause>] \perp]
Group by Clause :-	Group by [Cube Rollup] <Attribute Clause>
Having Clause :-	[Alias Aggregate Function Name (Attribute Name)] <Comparison Operator> [Attribute Value Attribute Value List]

Figure 6 – Modèle de requêtes de DWEB

Définissons la sémantique de la terminologie employée dans la figure 6.

- Les crochets [et] sont des délimiteurs.
- !<A> : A est requis.
- *<A> : A est optionnel.
- <A || B> : A ou B.
- <A | B> : A ou exclusif B.
- \perp : clause vide.
- Les éléments du langage SQL sont indiqués en gras.

4.2. Paramétrage

Comme les paramètres de la base de données de DWEB (section 3.2), ceux qui définissent sa charge (tableau 3) ont été conçus pour répondre au critère de simplicité

de Gray. Ils déterminent la manière dont le modèle de requêtes de la figure 6 s'instancie. Nous ne disposons ici que d'un nombre limité de paramètres de haut niveau (huit paramètres, puisque *PROB_EXTRACT* et *PROB_ROLLUP* sont calculés à partir de *PROB_OLAP* et *PROB_CUBE*, respectivement). Il n'est en effet pas envisageable d'entrer plus dans le détail si la taille de la charge atteint cinq cents requêtes comme c'est le cas dans TPC-DS, par exemple.

Nom du paramètre	Signification	Valeur par défaut
<i>NB_Q</i>	Nombre approximatif de requêtes dans la charge	100
<i>AVG_NB_ATT</i>	Nombre moyen d'attributs sélectionnés dans une requête	5
<i>AVG_NB_RESTR</i>	Nombre moyen de restrictions dans une requête	3
<i>PROB_OLAP</i>	Probabilité que le type de requête soit OLAP	0,9
<i>PROB_EXTRACT</i>	Probabilité que la requête soit une requête d'extraction	1 $-PROB_OLAP$
<i>AVG_NB_AGGREG</i>	Nombre moyen d'agrégations dans une requête OLAP	3
<i>PROB_CUBE</i>	Probabilité qu'une requête OLAP utilise l'opérateur <i>Cube</i>	0,3
<i>PROB_ROLLUP</i>	Probabilité qu'une requête OLAP utilise l'opérateur <i>Rollup</i>	1 $-PROB_CUBE$
<i>PROB_HAVING</i>	Probabilité qu'une requête OLAP contienne une clause <i>Having</i>	0,2
<i>AVG_NB_DD</i>	Nombre moyen de forages (<i>drill downs</i>) après une requête OLAP	3

Tableau 3 – Paramètres de la charge de DWEB

Remarque : *NB_Q* n'est qu'une approximation du nombre de requêtes parce que le nombre de forages (*drill downs*) effectués après une requête OLAP est variable. Nous ne pouvons donc arrêter le processus de génération que lorsque nous avons effectivement produit au moins *NB_Q* requêtes.

4.3. Algorithme de génération

L'algorithme de génération de la charge de DWEB est présenté dans les figures 7 et 8. Son objectif est de générer un ensemble de requêtes SQL-99 qui puisse être directement exécuté sur l'entrepôt de données synthétique défini à la section 3. Il est subdivisé en deux étapes :

1) génération d'une requête initiale qui peut être soit une requête OLAP, soit une requête d'extraction ;

2) si la requête initiale était de type OLAP, exécution d'un certain nombre de forages basés sur cette première requête. Plus précisément, à chaque fois qu'un nouveau forage est exécuté, un attribut d'un niveau inférieur de la hiérarchie d'une dimension est ajouté à la clause *Attribute Clause* de la requête précédente.

La première étape est elle-même subdivisée en trois sous-étapes :

- 1) les clauses *Select*, *From* et *Where* d'une requête sont générées simultanément en sélectionnant aléatoirement une table de faits et des dimensions (y compris un niveau hiérarchique dans chacune de ces dimensions) ;
- 2) la clause *Where* est complétée par des conditions supplémentaires ;
- 3) finalement, le choix du type de requête est effectué (requête OLAP ou requête d'extraction). Dans le second cas, la requête est terminée. Dans le premier, des fonctions d'agrégat appliquées à des mesures de la table de faits sont ajoutées à la requête, ainsi qu'une clause *Group by* pouvant inclure les opérateurs *Cube* ou *Rollup*. Il est également possible d'ajouter une clause *Having*. La fonction d'agrégat que nous appliquons sur les mesures est toujours une somme car c'est l'opération la plus courantes dans les cubes. De plus, les autres fonctions d'agrégat ont des complexités similaires en temps de calcul. Elles n'apporteraient donc pas plus d'information dans le cadre d'une étude de performance.

Nous utilisons trois classes de fonctions et une procédure dans cet algorithme.

1) `Random_string()` et `Random_float()` sont les mêmes fonctions que celles décrites dans la section 3.3. Cependant, nous introduisons ici la possibilité pour `Random_float()` d'utiliser une distribution aléatoire soit uniforme, soit gaussienne. Cela dépend des paramètres passés à la fonction : distribution uniforme pour un intervalle de valeurs, distribution gaussienne pour une valeur moyenne. Finalement, nous introduisons la fonction `Random_int()` qui se comporte exactement comme `Random_float()`, mais qui retourne des valeurs entières.

2) `Random_FT()` et `Random_dimension()` permettent de sélectionner une table de faits et une dimension qui décrit une table de faits donnée, respectivement. Elles utilisent toutes deux une distribution aléatoire gaussienne. `Random_dimension()` est également déjà décrite dans la section 3.3.

3) `Random_attribute()` et `Random_measure()` ont des comportements très similaires. Elles retournent un attribut ou une mesure, respectivement, appartenant à l'intention d'une table ou à une liste d'attributs. Elles utilisent toutes deux une distribution aléatoire gaussienne.

4) `Gen_query()` est la procédure qui est chargée de générer effectivement le code SQL-99 des requêtes de la charge, d'après les paramètres nécessaires pour instancier notre modèle de requêtes.

```

n = 0
While n < NB_Q do
  // Etape 1 : Requête initiale
  // Etape 1.2 : Clauses Select, From et Where
  i = Random_FT() // Sélection de la table de faits
  attribute_list = ∅
  table_list = ft(i)
  condition_list = ∅
  For k = 1 to Random_int(AVG_NB_ATT) do
    j = Random_dimension(ft(i)) // Sélection d'une dimension
    l = Random_int(1, ft(i).dim(j).nb_levels)
    // Positionnement au niveau hiérarchique l
    hl = ft(i).dim(j) // Niveau hiérarchique courant
    m = 1 // Compteur de niveau
    fk = ft(i).intention.primary_key.element(j)
    // (Cette clé étrangère correspond à ft(i).dim(j).primary_key)
    While m < l and hl.child ≠ NIL do
      // Construction de la jointure
      table_list = table_list ∪ hl
      condition_list = condition_list
        ∪ (fk = hl.intention.primary_key)
      // Niveau suivant
      fk = hl.intention.foreign_key
      m = m + 1
      hl = hl.child
    End while
    attribute_list = attribute_list
      ∪ Random_attribute(hl.intention)
  End for
  // Etape 1.2 : Compléter la clause Where
  For k = 1 to Random_int(AVG_NB_RESTRE) do
    condition_list = condition_list
      ∪ (Random_attribute(attribute_list) = Random_string())
  End for
  // Etape 1.3 : Sélection requête OLAP ou requête d'extraction
  p1 = Random_float(0,1)
  If p1 ≤ PROB_OLAP then // Requête OLAP
    // Agrégat
    aggregate_list = ∅
    For k = 1 to Random_int(AVG_NB_AGGREG) do
      aggregate_list = aggregate_list
        ∪ Random_measure(ft(i).intention)
    End for
  End if
n = n + 1
End while
../.

```

Figure 7 – Algorithme de génération de la charge – Partie 1

```

../..
    // Clause Group by
    group_by_list = attribute_list
    p2 = Random_float(0,1)
    If p2 ≤ PROB_CUBE then
        group_by_operator = CUBE
    Else
        group_by_operator = ROLLUP
    End if
    // Clause Having
    p3 = Random_float(0,1)
    If p3 ≤ PROB_HAVING then
        having_clause = (Random_attribute(aggregate_list), ≥,
            Random_float())
    Else
        having_clause = ∅
    End if
Else // Requête d'extraction
    group_by_list = ∅
    group_by_operator = ∅
    having_clause = ∅
End if
    // Génération de la requête SQL
    Gen_query(attribute_list, aggregate_list, table_list,
        condition_list, group_by_list, group_by_operator,
        having_clause)
    n = n + 1
    // Etape 2 : Eventuelles requêtes DRILL DOWN
    If p1 ≤ PROB_OLAP then
        k = 0
        While k < Random_int(AVG_NB_DD) and hl.parent ≠ NIL do
            k = k + 1
            hl = hl.parent
            att = Random_attribute(hl.intention)
            attribute_list = attribute_list ∪ att
            group_by_list = group_by_list ∪ att
            Gen_query(attribute_list, aggregate_list, table_list,
                condition_list, group_by_list, group_by_operator,
                having_clause)
        End while
        n = n + k
    End if
End while

```

Figure 8 – Algorithme de génération de la charge – Partie 2

5. Implémentation de DWEB

DWEB est implémenté sous la forme d'une application Java. Nous avons choisi le langage Java pour satisfaire le critère de portabilité de Gray. La version actuelle de notre prototype permet de générer des schémas en flocon de neige. Les schémas en constellation ne sont pas encore implémentés. Par ailleurs, comme les paramètres de DWEB peuvent sembler abstraits, notre prototype propose une estimation de la taille de l'entrepôt généré en méga-octets une fois que les paramètres sont fixés et avant la génération de la base de données. Ainsi, les utilisateurs peuvent ajuster les paramètres pour obtenir le type d'entrepôt de données qu'ils souhaitent. Notre prototype peut être interfacé à la plupart des systèmes de gestion de bases de données relationnels grâce à JDBC. La charge de DWEB est également en cours d'implémentation. Seule une charge simplifiée est disponible actuellement.

Par ailleurs, comme nous utilisons de nombreuses fonctions aléatoires, nous pensons aussi introduire dans notre prototype un générateur pseudo-aléatoire plus performant que ceux fournis en standard, comme celui proposé par Lewis et Payne [LEW 73], qui est un des meilleurs grâce à sa très grande période.

Finalement, bien que notre application soit en perpétuelle évolution, sa dernière version est en permanence disponible en ligne [DUC 04].

6. Conclusion et perspectives

Nous avons présenté dans cet article les spécifications complètes d'un nouveau banc d'essais pour entrepôts de données baptisé DWEB. La caractéristique principale de DWEB est qu'il permet de générer des entrepôts de données synthétiques variés, ainsi que les charges (ensemble de requêtes) associées. Les schémas d'entrepôt classiques, comme les modèles en étoile, en flocon de neige et en constellation sont en effet supportés. Nous considérons DWEB comme un banc d'essais d'ingénierie destiné aux concepteurs d'entrepôts de données et de systèmes. En cela, il n'est pas concurrent du banc d'essais TPC-DS (actuellement en cours de développement), qui est plus à destination des utilisateurs finaux, pour des évaluations de performances « pures ». Il faut toutefois remarquer que le schéma d'entrepôt de données de TPC-DS peut être modélisé à l'aide de DWEB. Cependant, la charge de DWEB n'est actuellement pas aussi élaborée que celle de TPC-DS. Finalement, nous avons pour objectif de fournir dans cet article les spécifications les plus complètes possibles de DWEB, de manière à ce que notre banc d'essais puisse être implémenté facilement par d'autres chercheurs et/ou concepteurs d'entrepôts de données.

Nos travaux futurs dans ce domaine sont divisés en quatre axes. Premièrement, il nous est nécessaire de terminer l'implémentation de DWEB (charge, schémas en constellation avec hiérarchies partagées, génération de fichier plats et chargement dans une base de données, tests divers...). Ce travail est actuellement en cours. Des expériences avec DWEB devraient également nous permettre de proposer un meilleur paramétrage par défaut de notre banc d'essais. Nous encourageons également d'autres

chercheurs et/ou concepteurs d'entrepôts de données à publier les résultats de leurs propres expériences avec DWEB.

Deuxièmement, ils nous faut tester la pertinence (d'après la définition de Gray) de notre banc d'essais pour l'évaluation de performance des entrepôts de données dans un contexte d'ingénierie. Afin d'atteindre cet objectif, nous pensons comparer l'efficacité de plusieurs techniques de sélection automatique d'index et de vues matérialisées (dont certaines de nos propositions, ce qui a constitué une des motivations de la conception de DWEB à l'origine).

Troisièmement, nous comptons progressivement améliorer DWEB. Par exemple, dans cet article, nous supposons implicitement que nous pouvons réutiliser ou facilement adapter le protocole d'exécution et les mesures de performance de TPC-DS. Une réflexion plus élaborée à propos de ces sujets pourrait être intéressante. Il sera également important d'inclure complètement le processus d'ETL dans notre charge. Plusieurs travaux existants [LAB 98] pourraient d'ailleurs nous y aider. Finalement, il serait intéressant de rendre notre charge plus dynamique, comme cela a été fait dans la plate-forme d'évaluation de performance DoEF [HE 03, HE 04]. Bien que la dynamique ne soit peut-être pas un facteur pertinent dans une charge décisionnelle, cette possibilité pourrait néanmoins être explorée. Pour finir, puisque nous travaillons par ailleurs sur l'entreposage de données complexes (multiformats, multistruktures, multisources, multimodales et/ou multiversions), nous envisageons également une extension « données complexes » pour DWEB à long terme.

Finalement, DWEB pourrait évoluer pour devenir capable de proposer automatiquement des configurations d'entrepôts de données et des charges afin d'évaluer les performances d'une architecture d'entrepôt donnée ou d'une technique de sélection d'index ou de vue matérialisée. DWEB pourrait alors capitaliser la connaissance acquise au cours de diverses expériences et la réutiliser pour proposer de nouvelles configurations, plus pertinentes.

7. Bibliographie

- [BAL 93] BALLINGER C., « *TPC-D : Benchmarking for Decision Support* », The Benchmark Handbook for Database and Transaction Processing Systems, Morgan Kaufmann, 1993.
- [BHA 96] BHASHYAM R., « TCP-D : The Challenges, Issues and Results », *SIGMOD Record*, vol. 25, n° 4, 1996, p. 89-93.
- [BMC 00] BMC SOFTWARE, « Performance Management of a Data Warehouse », <http://www.bmc.com>, 2000.
- [DAR 00] DARMONT J., SCHNEIDER M., « Benchmarking OODBs with a Generic Tool », *Journal of Database Management*, vol. 11, n° 3, 2000, p. 16-27.
- [DEM 95] DEMAREST M., « A Data Warehouse Evaluation Model », *Oracle Technical Journal*, vol. 1, n° 1, 1995, page 29.
- [DUC 04] DUCREUX S., PENIN P.-M., « DWEB Java prototype », <http://bdd.univ-lyon2.fr/download/dweb.tgz>, 2004.

- [GRA 93] GRAY J., *The Benchmark Handbook for Database and Transaction Processing Systems*, Morgan Kaufmann, second édition, 1993.
- [GRE 04a] GREENFIELD L., « Performing Data Warehouse Software Evaluations », <http://www.dwinfocenter.org/evals.html>, 2004.
- [GRE 04b] GREENFIELD L., « What to Learn About in Order to Speed Up Data Warehouse Querying », <http://www.dwinfocenter.org/fstquery.html>, 2004.
- [HE 03] HE Z., DARMONT J., « DOEF : A Dynamic Object Evaluation Framework », *14th International Conference on Database and Expert Systems Applications (DEXA 03)*, Prague, Czech Republic, vol. 2736 de *LNCIS*, September 2003, p. 662-671.
- [HE 04] HE Z., DARMONT J., « Une plate-forme dynamique pour l'évaluation des performances des bases de données à objets », *Ingénierie des Systèmes d'Information*, vol. 9, n° 1, 2004, p. 109-127.
- [INM 02] INMON W., *Building the Data Warehouse*, John Wiley & Sons, third édition, 2002.
- [KIM 02] KIMBALL R., ROSS M., *The Data Warehouse Toolkit : The Complete Guide to Dimensional Modeling*, John Wiley & Sons, second édition, 2002.
- [LAB 98] LABRINIDIS A., ROUSSOPOULOS N., « A Performance Evaluation of Online Warehouse Update Algorithms », rapport n° CS-TR-3954, November 1998, Department of Computer Science, University of Maryland.
- [LEW 73] LEWIS T., PAYNE W., « Generalized feedback shift register pseudorandom number algorithm », *ACM Journal*, vol. 20, n° 3, 1973, p. 458-468.
- [Obj03] Object Management Group, « Common Warehouse Metamodel (CWM) Specification version 1.1 », March 2003.
- [PEN 03] PENDSE N., « The OLAP Report : How not to buy an OLAP product », http://www.olapreport.com/How_not_to_buy.htm, December 2003.
- [POE 00] POESS M., FLOYD C., « New TPC Benchmarks for Decision Support and Web Commerce », *SIGMOD Record*, vol. 29, n° 4, 2000, p. 64-71.
- [POE 02] POESS M., SMITH B., KOLLAR L., LARSON P.-A., « TPC-DS : Taking Decision Support Benchmarking to the Next Level », *ACM SIGMOD 2002, Madison, USA*, June 2002.
- [POO 03] POOLE J., CHANG D., TOLBERT D., MELLOR D., *Common Warehouse Metamodel Developer's Guide*, John Wiley & Sons, 2003.
- [Tra98] Transaction Processing Performance Council, « TPC Benchmark D Standard Specification version 2.1 », February 1998.
- [Tra03a] Transaction Processing Performance Council, « TPC Benchmark H Standard Specification version 2.1.0 », August 2003.
- [Tra03b] Transaction Processing Performance Council, « TPC Benchmark R Standard Specification version 2.1.0 », August 2003.